

Multi-Sourced Data Retrieval in Groomed Elastic Optical Networks

Juzi Zhao and Vinod M. Vokkarane

Department of Electrical and Computer Engineering
University of Massachusetts Lowell, MA 01854
Email: {juzi_zhao, vinod_vokkarane}@uml.edu

Abstract—Erasure coding has been widely adopted by data center networks, where the data is encoded and stored in multiple locations. Therefore, an efficient data retrieval service is needed to transfer encoded data from replicated stored nodes to a single destination. Elastic Optical Networks are a promising backbone technology for data center communication due to their capability to efficiently and flexibly allocate the huge optical bandwidth to heterogeneous traffic demands. Traffic grooming is a technique to increase the spectrum efficiency and transponder utilization of optical networks. In this paper, the traffic-grooming-enabled erasure-coded multi-sourced data retrieval routing and scheduling problem is studied for dynamic traffic in elastic optical networks. A novel heuristic is proposed. Numerical results indicate that traffic grooming decreases the blocking ratio by a factor of 3 and cost by 22% compared with a baseline algorithm.

Index Terms—Traffic grooming, elastic optical networks, data retrieval, spectrum allocation, dynamic traffic.

I. INTRODUCTION

Taking of the advantage of its intrinsic flexibility and high efficiency in allocating the optical spectrum resources, *Elastic Optical Networks (EON)* are an ideal candidate for next generation data center networks to support the explosively increasing demands of highly-diversified traffic from various bandwidth-hungry applications, such as e-science, cloud and grid computing, and inter data center communications, which require large volumes of datasets to be distributed and processed by geographically disperse users. Optical Orthogonal Frequency Division Multiplexing (O-OFDM), Bandwidth-Variable Cross-Connect (BV-OXC), and Bandwidth-Variable Transponder (BV-T) are important technologies to enable elastic optical networks. The optical spectrum of each fiber (e.g., the C-band) is divided into finer granularity subcarriers (e.g., 6.25 GHz or 12.5 GHz) than the bandwidth of the wavelengths in traditional rigid wavelength division multiplexing (WDM)-based optical network (e.g., 50 GHz). A variety of modulation schemes (e.g.,

quadrature phase-shift keying (QPSK) and 16-point quadrature amplitude modulation (16-QAM)) can be adopted by different subcarriers and result in different bit rates. Therefore, only the required number of *contiguous* subcarriers needs to be allocated to a service in order to satisfy its throughput requirement [1]. The physical impairments in elastic optical networks are commonly represented by a transmission reach limit for each particular modulation scheme. In addition, to prevent the interference among connections that share same link at any time throughout their durations, a guardband is placed between two adjacent subcarrier bands assigned to different connections [1]–[3].

Traffic grooming is a widely-taken approach to increase the bandwidth efficiency, device utilization, energy consumption, and network cost of elastic optical networks. By aggregating several low-speed traffic connections into a single high-speed lightpath by electrical switching, traffic grooming can reduce the spectrum resources by allocating consecutive subcarriers without interim guardbands; in addition, only one single bandwidth variable transponders (BVT) is needed to accommodate the traffic. The authors of [4] introduced traffic grooming to EON and showed that traffic grooming can lead to significant savings in transponder and spectrum allocation. A new traffic grooming algorithm for the connection establishment of deadline-driven requests is proposed in [8], which grooms batches of requests to establish lightpaths with diverse bandwidth demands with deadline requirements. The authors of [6] propose a three-layered auxiliary graph (AG) model to address mixed-electrical-optical grooming under dynamic traffic scenario. Traffic grooming has also been successfully applied to multipath routing algorithms to increase bandwidth blocking ratio and BVT usage [7], [8].

In data centers, erasure coding has been widely adopted to protect against disk and node failures to increase reliability [9], such as Windows Azure [10], Facebook [11] and Google [12]. Under a (n, m) erasure coding, data is encoded and stored in n storage nodes

such that the pieces stored in any m of these n nodes suffice to recover the entire data. It is desired to provide fast and convenient erasure-coded data retrieval from distributed repositories to a single site. A user requesting erasure-coded data retrieval from n remote sites (each has a piece of the erasure coded data) can have a choice of m selected storage sites between the n total storage sites. The speed to transfer large datasets can be efficiently improved by parallel transmission through multiple paths from one selected storage site to the destination.

A related problem studied by researchers is the many-to-one data aggregation, where a set of data, located at distributed databases in the network, need to be transmitted to a single destination. There are several research papers about dynamic many-to-one data aggregation service in WDM optical networks. A hybrid approach that combines offline and online scheduling for the problem was proposed in [13]. The authors of [14] propose four dynamic scheduling algorithms to address the issue of large file transfers on multi-user optical grid network. Algorithms for routing, wavelength assignments, grooming, and scheduling the files are proposed in [15] with the assumption that wavelength converter are placed at every node in the network (opaque networks). Scheduling and resource allocation algorithms are proposed in [16], [17] for divisible load applications in multidomain optical grid. However, none of them considered erasure-coded data aggregation, which includes storage node selection. Our previous paper [18] investigates the erasure-coded data retrieval problem for dynamic traffic in elastic optical networks without traffic grooming and transponder capacity limit.

To the best of our knowledge, this is the first work that study the dynamic erasure-coded traffic-grooming-enabled data retrieval from multiple repository sources in elastic optical networks with the consideration of traffic grooming, storage node selection, physical layer impairments, multi-path routing, different modulation assignment, and contiguous subcarriers allocation.

The paper is organized as follows. In Section II, we present the problem statement. The proposed heuristic and benchmark algorithm are presented in Section III. Section IV presents and discusses the numerical results. We conclude the paper in Section V.

II. PROBLEM STATEMENT

The network is modeled as a graph $G(V, E)$ with node set V and link set E . S subcarriers are available on each link $e \in E$, and the bandwidth of each subcarrier is C GHz. Each node $v \in V$ is equipped with

T bandwidth variable transponders for launching and receiving lightpaths; these transponders are shared among all lightpaths using that using node v as an ending node (share-by-node model), and the capacity of each transponder is Q Gbps. Each lightpath can be assigned any of B modulation formats. There is a transmission reach limit of each modulation format due to physical layer impairments; thus the distance of a lightpath assigned modulation format $b \in B$ should be no greater than the reach limit R_b . We assume that there are K pre-computed shortest (in distance) paths for each pair of nodes in the network. Based on the distance of each path and the transmission reach limits of the modulation formats, the modulation of each path can be obtained. For a lightpath with modulation b , if w subcarriers on it are allocated, then the required transponder resource (in bit rate) at the two ending nodes of the lightpath is $M_b w C$ Gbps, where M_b is the spectrum efficiency of modulation scheme b (in terms of Gbps/GHz). g subcarriers served as guardband should be placed between two adjacent subcarrier bands assigned to different lightpaths which share same link at any time throughout their durations to avoid interference at intermediate switches. We assume a slotted time system, and any unit of work that alters the state of the network must occur at the beginning of a time slot and must finish at an integer multiple of the time slot span.

For dynamic traffic, erasure-coded data retrieval requests randomly arrive to the network. Each request with erasure code (n, m) specifies its destination node d and the n storage nodes. There is a data piece on each storage node, resulting in n different data pieces in total. To retrieve the entire data set, any m of n data pieces need to be transmitted to the destination. Suppose the minimum granularity of data size is F , i.e., each data piece is an integer multiple of F , and we call the minimum data granularity a *data segment*. Denote H_i as the number of data segments of each data piece of request i , thus, the total number of required data segments to retrieve the whole data set of request i is mH_i . Different data segments of one data piece could be transferred to the destination through multiple paths. There is also a deadline requirement D associated with each request, if the entire data retrieval cannot finish within D time slots, the request is blocked. The objective is to minimize the request blocking ratio, which is defined as the number of blocked requests over the total number of requests arrived to the network, with single-hop and multi-hop traffic grooming.

For each accepted request, the m (out of n) storage

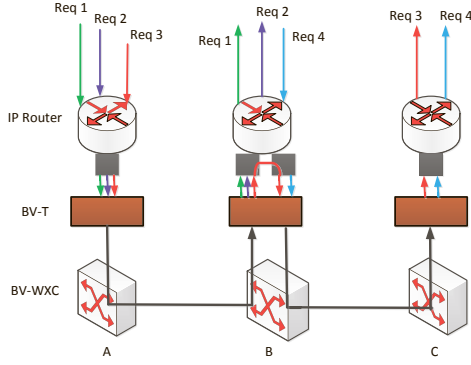


Fig. 1: The fundamental architecture of an EON.

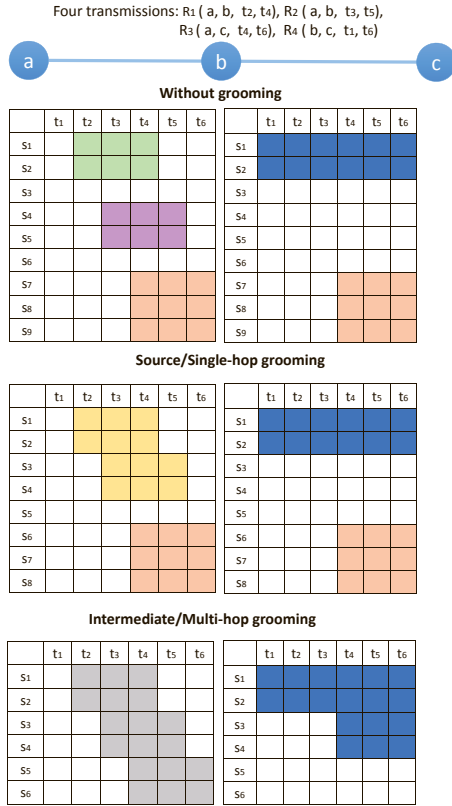


Fig. 2: An example of traffic grooming.

nodes need to be selected, in turn, a route, a modulation format, a *contiguous* subcarrier band, a starting transmitting time slot, and transponder resources at the storage node and destination node need to be allocated to each data segment on the m storage nodes to satisfy the deadline requirement of an accepted request.

Multiple low-speed traffic connections can be electronically groomed into a single connection and dropped to/added from the optical layer by a single router port as shown in Fig. 1. A traffic connection for a data

segment transmission can be groomed with other data segment transmissions at its source node (i.e., single hop grooming) and/or an intermediate node (i.e., multi-hop grooming). A transiting connection is retransmitted by transponders at the intermediate grooming nodes. Besides increased subcarrier and transponder resource utilization, other advantages of traffic grooming in EONs include transmission reach extension, modulation format conversion, and subcarrier-band conversion. In the example Fig. 2, suppose there exists one data segment transmission R_4 from Node b to Node c during the time period t_1-t_6 . A lightpath has been established for it with two subcarriers (s_1, s_2) allocation on Link ($b-c$). Also suppose that there are three data segment transmissions, two from Node a to Node b , and one from Node a to Node c . The first transmission R_1 takes three time slots, t_2-t_4 , the second transmission R_2 takes time slots t_3-t_5 , while the third transmission R_3 takes time slots t_4-t_6 . Suppose the path ($a-b-c$) is allocated to all of them. Suppose that according to the transmission reach limits with different modulation level formats, two subcarriers are enough to transmit each data segment on Links ($a-b$) and ($b-c$) by adopting a higher modulation format than the one for path ($a-b-c$). In the case that there is no grooming, one transponder at Nodes a and b (say, transponder 1) needs to be allocated for R_1 , another transponder (e.g., transponder 2) is required to be allocated to R_2 , and a transponder at Nodes a and c (say, transponder 3) needs to be allocated for R_3 . In addition, with one subcarrier as guardband, subcarriers s_1, s_2 are allocated to R_1 , subcarriers s_4, s_5 are allocated to R_2 , and subcarriers s_7, s_8, s_9 are allocated to R_3 (where subcarriers s_3 and s_6 serve as guardbands). In the case that only source grooming/single-hop grooming is allowed, the two transmissions R_1 and R_2 can be groomed together. Only one transponder at Nodes a and b is allocated to the newly established lightpath for these two requests. Furthermore, there is no guardband needed between the subcarrier allocation for the two transmissions, in turn, only 8 subcarriers are allocated in total. In the case that intermediate grooming/multi-hop grooming is also applied, more subcarriers can be saved as shown in the example by grooming the R_3 into the existing lightpath from Node b to Node c . The intermediate grooming Nodes b includes Optical-Electrical-Optical (O-E-O) conversion, thus, they can work as spectrum converters and modulation format converters. One possible subcarrier allocation is shown in the example.

III. PROPOSED ALGORITHMS AND BENCHMARK

A. Grooming-enabled Minimum Resource Algorithm (GMinR)

We proposed a heuristic, GMinR, to solve this problem, which takes in to account traffic grooming at intermediate nodes and at source node when scheduling resources (transponders and subcarriers) to the data segments of the newly arrived request i with arrival time A_i , deadline D_i , and destination node d . The algorithm for each incoming request i is summarized in Algorithm 1.

The first step of GMinR is to select the m of n storage nodes (Lines 1-20). Each candidate storage node v is assigned a weight $\gamma_v = \alpha \sum_{t=A_i}^{A_i+D_i} \tau_v^t + (1-\alpha)\sigma_v$, wherein τ_v^t is the available transponder capacity of node v at time slot t , σ_v is the available subcarriers capacity on the K paths from node v to destination node d during time period from A_i to $A_i + D_i$, b_k is the modulation scheme with the k^{th} path, ϕ_e^{st} is a binary variable, its value is 1 if the subcarrier s on link e is available at time slot t , ϕ_k^{st} is also a binary variable, its value is 1 if the subcarrier s on the k^{th} path is available at time slot t (where means that $\phi_e^{st} = 1$ for all e on the k^{th} path), and α is an input parameter to balance the two kinds of resources. Then, the m storage nodes with larger weights are selected. The idea behind is to select storage nodes with many available resources for the request, thereby leaving more resources for future requests. Then the selected mH_i data segments are scheduled one by one as follows.

For each data segment $1 \leq j \leq mH_i$ (Lines 22-30), a logical graph $G'_{t_1, t_2}(V', E')$ is created for every time duration t_1 to t_2 ($A_i \leq t_1 \leq t_2 \leq A_i + D_i$), wherein t_1 and t_2 are starting transmission time slot and finishing transmission time slot of data segment j respectively. Each existing lightpath is checked whether having enough remaining transponder capacity and subcarriers (more available subcarriers may need to be allocated to this existing lightpath) to transfer data segment j , and let these lightpaths be set L . The logical nodes V' includes the destination node d , the storage node of data segment j , and the two ending nodes of each existing lightpath $l \in L$. A logical link $e' \in E'$ exists between each pair of logical nodes v'_a and v'_b . Let P denote a set including both existing lightpaths ($\in L$) between v'_a and v'_b and the K precomputed shortest distance paths for node pair

v'_a and v'_b . Each logical link is assigned a weight by equation (1), wherein $p \in P$, w is the number of allocated subcarriers to accommodate data segment j (for existing lightpath $p \in L$, $w \geq 0$ are the additional subcarriers allocated to p in order to accommodate j), $\eta_{v'_a}^{pw}$ and $\eta_{v'_b}^{pw}$ are the required transponder resource at node v'_a and v'_b , respectively, which are calculated as $M_b(p)wC$, where $M_b(p)$ is the spectrum efficiency for the modulation scheme with p , and N_p is the number of hops on p . Dijkstra minimum weight algorithm is used for the logical graph to get the end-to-end route r for data segment j . The weight of the route is $\beta_{G'}^{t_1, t_2} = \sum_{e' \in r} \beta_{e'}^{t_1, t_2}$. At last, we can obtain the weight of each data segment j as $\omega_j = \min_{t_1, t_2} \beta_{G'}^{t_1, t_2}$.

After obtaining the weights of all the mH_i data segments, the data segment with minimum weight $\min_j \omega_j$ is selected, and moved to set J , the resources allocation achieving the weight $\min_j \omega_j$ are also kept into records. The weights of remaining data segments (not in set J) are re-calculated under the condition that the resource allocation for data segments in set J are reserved. This procedure is repeated until all data segments are in set J . If one of the mH_i data segments cannot be allocated within the deadline D_i , the whole request is blocked. Otherwise, the request is accepted, corresponding lightpaths are established or updated (if they are existing lightpaths), and the corresponding resources in the network are scheduled. The subcarrier allocation method is First Fit, wherein the first subcarrier band with the required available subcarriers is allocated.

B. Source Grooming-enabled Minimum Resource Algorithm (SGMinR).

The SGMinR heuristic is the case only source grooming is applied. The difference compared to GMinR algorithm is that the existing lightpaths which are considered when creating each logical graph $G'_{t_1, t_2}(V', E')$ should have the same source node as the data segment under consideration.

C. Baseline

The baseline heuristic is the case without traffic grooming. The differences compared to GMinR algorithm are: a) the logical nodes in the logical graph $G'_{t_1, t_2}(V', E')$ are the source node v and destination node d ; b) only the K precomputed shortest distance paths for node pair v, d are included in the set P .

$$\beta_{e'}^{t_1, t_2} = \min_{p \in P} \alpha \frac{2(\eta_{v'_a}^{pw} + \eta_{v'_b}^{pw})(t_2 - t_1 + 1)}{Q} + (1 - \alpha) \frac{N_p w (t_2 - t_1 + 1)}{S}, \quad (1)$$

```

1 Set  $\mu_e^{st} = 0, \sigma_v = 0 \forall v, s, t, e$ 
2 foreach node  $v$  of  $n$  candidate nodes do
3   Find  $\tau_v^t = \min(T_v^t, T_d^t)$ 
4   Find  $\phi_k^{st}$  state based on  $\phi_e^{st}$  for all
    $1 \leq k \leq K, e \in p, 1 \leq s \leq S, A_i \leq t \leq A_i + D_i$ 
5   foreach  $1 \leq s \leq S$  do
6     foreach  $A_i \leq t \leq A_i + D_i$  do
7       foreach  $1 \leq k \leq K$  do
8         Find  $\mu_k^{st} = \sum_{e \in p} \mu_e^{st}$ 
9         if  $\phi_k^{st} = 1$  and  $b_k > \mu_k^{st}$  then
10          foreach link  $e$  on path  $k$  do
11             $\mu_e^{st} = b_k$ 
12          end
13           $\sigma_v = \sigma_v + b_k - \mu_k^{st}$ 
14        endif
15      end
16    end
17  end
18   $\gamma_v = \alpha \sum_{t=A_i}^{A_i+D_i} \tau_v^t + (1 - \alpha)\sigma_v$ 
19 end
20 Sort the candidate nodes by decreasing order
  of  $\gamma_v$ , and select the first  $m$  source nodes
21 while there are still unallocated data segments
  do
22   foreach node  $v$  of  $m$  candidate nodes do
23     foreach data segment  $j$  of request  $i$  on
     the storage node  $v$  do
24       foreach  $t_1, t_2$  do
25         Create logical graph  $G'_{t_1, t_2}(V', E')$ 
26         Find the weight of each logical
         link  $e' \in E'$  by equation (1)
27         Use Dijkstra minimum weight
         algorithm to get the weight of
         route  $\beta_{G'}^{t_1, t_2}$ 
28       end
29        $\omega_{vj} = \min_{t_1, t_2} \beta_{G'}^{t_1, t_2}$ 
30     end
31   end
32   Select the  $j, v$  pair that achieves the
   minimum  $\omega_{vj}$ 
33   Reserve the corresponding resources
34 end
35 if all the data segments on all selected  $m$ 
  storage nodes are successfully allocated then
36   Request  $i$  is accepted, schedule resources
37 endif
38 else
39   Request  $i$  is blocked
40 endif

```

Algorithm 1: GMinR algorithm

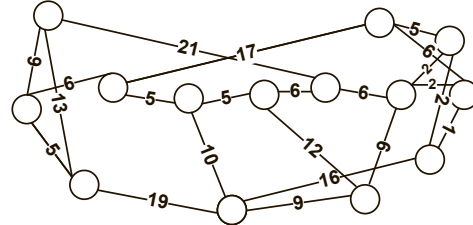


Fig. 3: NSFNET topology.

TABLE I: Parameters

S, Number of frequency-slots per link	320
Spectrum efficiency of BPSK	1 bps/Hz
Spectrum efficiency of QPSK	2 bps/Hz
Spectrum efficiency of 8QAM	3 bps/Hz
Spectrum efficiency of 16QAM	4 bps/Hz
C, bandwidth of a frequency per link	12.5 Gbps
g , number of slots for guard-band per request	1
Transmission reach of BPSK	5000 km
Transmission reach of QPSK	2500 km
Transmission reach of 8-QAM	1250 km
Transmission reach of 16-QAM	625 km

IV. NUMERICAL RESULTS

We present simulation results for the 14-node NSF network topology in Fig. 3. The number on each link corresponds to the link length in units of 100 km. 4 modulation formats: BPSK, QPSK, 8-QAM, and 16-QAM can be assigned to lightpaths. The parameters related to physical impairments are listed in Table I [3]. There are $K = 3$ precomputed shortest distance paths for each node pair [19]. Reed-Solomon code (9, 6) in GFS II in Google [12] is adopted. The capacity of each transponder is 400Gbps. For each simulation, 10000 requests are generated following a Poisson process with arrival rate. The destination node and 9 storage sources nodes of each request is uniformly selected from all the nodes in the network. The number of data segments per piece, H , for each request is randomly selected in the range (3, 5). Parameter α is set as 0.5 to balance the allocation of transponders and subcarrier resources. Each result point is an average value over ten simulation seeds.

The resulting request blocking ratio as a function of the number of transponders at each node is shown in Fig. 4. The deadline for each request is set as 50 time slots and the average request arrival rate is 30 new requests per time slots. It is shown that the proposed GMinR and SGMinR reduces the request blocking ratio comparing to the baseline without traffic grooming. SGMinR with 60 transponders per node has better blocking performance than baseline algorithm with 65 transponders per node. Furthermore, GMinR with 60

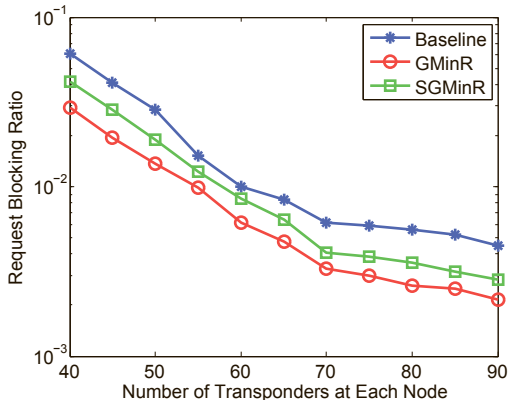


Fig. 4: Request blocking vs. number of transponders per node.

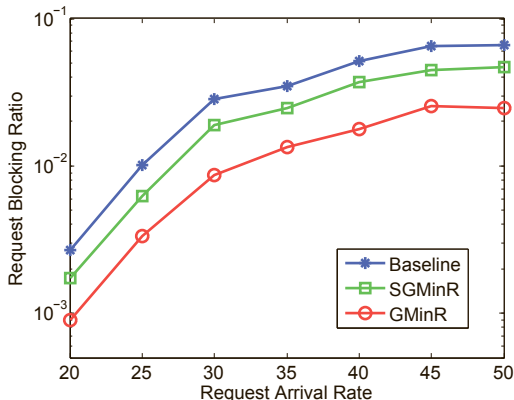


Fig. 5: Request blocking vs. request arrival rate.

and 70 transponders per node have better blocking performance than baseline algorithm with 70 and 90 transponders per node, respectively, indicating that GMinR could efficiently groom multiple data retrieval transmissions and save the transponder resources (up to 22%).

The multi-hop grooming needs O-E-O conversion, which results in delay (due to store-and-forwarding) and additional energy consumption, therefore, it is better to limit the number of intermediate groomings per data segment transmission. In the simulations, the maximum number of logical hops per data segment transmission is 4 (at most 3 intermediate groomings on one end-to-end path), and the average number of logical hops for each data segment is in the range 1.132-1.093, which slightly decreases as the number of transponders per node increases.

The resulting request blocking ratio as a function of the arrival rate is shown in Fig. 5. The deadline for each request is set as 50 time slots from arrival time slot and there are 50 transponders at each node. It

is shown that the proposed GMinR can significantly reduce the blocking ratio (up to a factor of 3 at load 30 Erlang) compared with the baseline heuristic. The corresponding results for number of logical hops per data segment transmission confirm our previous observation that there are not too many intermediate groomings involved in the data retrieval transmissions. The maximum number of logical hops per data segment transmission is still 4, and the average number of logical hops for each data segment is in the range 1.091-1.184, which slightly increases as the request arrival rate increases.

V. CONCLUSIONS

By including traffic-grooming, a novel routing and spectrum allocation algorithm was derived for the dynamic multi-sourced erasure-coded data retrieval problem in elastic optical networks. Simulations results showed that the blocking ratio can be significantly reduced (up to a factor of 3) in comparison with a baseline method, in addition, the number of transponders per node can be decreased 22% to achieve the same blocking performance. In the future work, we plan to investigate hop-constrained grooming policies, batch allocation, and traffic grooming with sliceable bandwidth variable transponders [20].

ACKNOWLEDGMENT

This work has been supported by the DOE PROPER project under grant DE-SC0012115TDD, NSF CAR-GONET project under grant CNS-1406370 and NSF CCDNI project under grant ACI-1541434.

REFERENCES

- [1] K. Christodoulopoulos, I Tomkos, and E. Varvarigos, "Elastic Bandwidth Allocation in Flexible OFDM-Based Optical Networks," *Journal of Lightwave Technology*, vol. 29, pp. 1354–1366, 2011.
- [2] X. Wan, N. Hua, and X. Zheng, "Dynamic Routing and Spectrum Assignment in Spectrum-Flexible Transparent Optical Networks," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 4, no. 8, pp. 603–613, 2012.
- [3] W. Lu and Z. Zhu, "Dynamic Service Provisioning of Advance Reservation Requests in Elastic Optical Networks," *Journal of Lightwave Technology*, vol. 31, no. 10, pp. 1621–1627, 2013.
- [4] G. Zhang, M. De Leenheer, and B. Mukherjee, "Optical Traffic Grooming in OFDM-based Elastic Optical Networks," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 4, no. 11, pp. B17B25, 2012.
- [5] P. Moura, N. da Fonseca, and R. Scaraficci, "Traffic Grooming of Batches of Deadline-driven Requests in Elastic Optical Networks", in *IEEE Global Communications Conference (GLOBE-COM)*, Austin, TX, Dec. 2014.
- [6] J. Zhang, Y. Ji, M. Song, X. Yu, J. Zhang, and B. Mukherjee, "Dynamic Traffic Grooming in Sliceable Bandwidth-variable Transponder-enabled Elastic Optical Networks," *Journal of Lightwave Technology*, vol. 33, no. 1, pp. 183–191, 2015.

- [7] Z. Fan, Y. Qiu, and C.-K. Chan, "Dynamic Multipath Routing with Traffic Grooming in OFDM-based Elastic Optical Path Networks," *Journal of Lightwave Technology*, vol. 33, no. 1, pp. 275–281, 2015.
- [8] M. N. Dharmaweera, J. Zhao, L. Yan, M. Karlsson, and E. Agrell, "Traffic-grooming-and Multipath-routing-enabled Impairment-aware Elastic Optical Networks," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 8, no. 2, pp. 58–70, 2016.
- [9] C. Angllano, R. Gaeta and M. Grangetto, "Exploiting Rateless Codes in Cloud Storage Systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 5, pp. 1313–1322, May 2015.
- [10] C. Huang, H. Simitci, Y. Xu, A. Ogus, B. Calder, P. Gopalan, J. Li, and S. Yekhanin, "Erasure Coding in Windows Azure Storage," *USENIX ATC*, 2012.
- [11] D. Borthakur et al., "HDFS RAID," *Hadoop User Group Meeting*, Nov. 2010.
- [12] D. Ford, F. Labelle, F.I. Popovici, M. Stokely, V. Truong, L. Barroso, C. Grimes, and S. Quinlan, "Availability in Globally Distributed Storage Systems," in *Proc. USENIX OSDI*, 2010.
- [13] A. Banerjee, Wu-chun Feng, D. Ghosal and B. Mukherjee, "Algorithms for Integrated Routing and Scheduling for Aggregating Data from Distributed Resources on a Lambda Grid", *IEEE Transactions on Parallel and Distributed Systems*, vol. 19, no.1, pp.24–34, Jan. 2008.
- [14] M. Hu, W. Guo, and W. Hu, "Dynamic Scheduling Algorithms for Large File Transfer on Multi-user Optical Grid Network Based on Efficiency and Fairness", in *Fifth International Conference on Networking and Services*, Valencia, Spain, Apr. 2009.
- [15] D. Andrei, M. Tornatore, D. Ghosal, C.U. Martel, and B. Mukherjee, "On-Demand Provisioning of Data-Aggregation Sessions Over WDM Optical Networks", *Journal of Lightwave Technology*, vol. 27, no. 12, pp. 1846 – 1855, June 2009.
- [16] M. Abouelela and M. El-Darieby, "Multidomain Hierarchical Resource Allocation for Grid Applications," *Journal of Electrical and Computer Engineering - Special issue on Resource Allocation in Communications and Computing*, January 2012.
- [17] M. Abouelela and M. El-Darieby, "Scheduling Big Data Applications Within Advance Reservation Framework in Optical Grids," *Applied Soft Computing*, vol. 38, pp. 1049–1059, Jan. 2016.
- [18] J. Zhao and V. M. Vokkarane, "Dynamic Erasure-coded Data Retrieval in Elastic Optical Data Center Networks," in *Proc. IEEE Sarnoff*, Sept. 2016.
- [19] J. Y. Yen, "Finding the K Shortest Loopless Paths in a Network," *Manage. Sci.*, vol. 17, no. 11, pp. 712–716, July 1971.
- [20] N. Sambo, P. Castoldi, A. D'Errico, E. Riccardi, A. Pagano, M. Moreolo, J. Fabrega, D. Rafique, A. Napoli, S. Frigerio, E. Salas, G. Zervas, M. Nolle, J. Fischer, A. Lord, J. Gimenez, "Next generation sliceable bandwidth variable transponders", *IEEE Communications Magazine*, vol. 53, no. 2, pp. 163–171, 2015.